

Sigmoid function $\sigma(\eta)$

$$\sigma(\eta) = \frac{1}{1 + e^{-\eta}} = \frac{e^{\eta}}{1 + e^{\eta}}, \quad 0 < \sigma(\eta) < 1$$

$$\frac{d\sigma(\eta)}{d\eta} = -\frac{-e^{-\eta}}{(1 + e^{-\eta})^2} = \frac{e^{-\eta}}{(1 + e^{-\eta})^2} = \frac{1}{1 + e^{-\eta}} \left(\frac{e^{-\eta}}{1 + e^{-\eta}} \right) = \frac{1}{1 + e^{-\eta}} \left(1 - \frac{1}{1 + e^{-\eta}} \right) = \sigma(\eta) [1 - \sigma(\eta)]$$

$$\frac{d \log \sigma(\eta)}{d\eta} = \frac{1}{\sigma(\eta)} \cdot \frac{d\sigma(\eta)}{d\eta} = 1 - \sigma(\eta)$$

Logistic Regression - Representation

Input: $\mathbf{x} \in \mathbb{R}^D$

Output: $y \in \{0, 1\}$

Training data: $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$

Model:

$$p(y = 1 | \mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$$

$$p(y = 0 | \mathbf{x}, \mathbf{w}) = 1 - \sigma(\mathbf{w}^T \mathbf{x}), \quad \sigma(\eta) = \frac{1}{1 + e^{-\eta}}$$

$$f(\mathbf{x}) : \mathbf{x} \rightarrow y, \quad f(\mathbf{x}) = \begin{cases} 1, & p(y = 1 | \mathbf{x}, \mathbf{w}) > 0.5 \\ 0, & \text{otherwise} \end{cases} = \begin{cases} 1, & \sigma(\mathbf{w}^T \mathbf{x}) > 0.5 \\ 0, & \text{otherwise} \end{cases}$$

Model parameters: Weights $\mathbf{w} \in \mathbb{R}^D$ (to be learned)

Logistic Regression - Evaluation Criterion

Data likelihood for 1 training sample:

$$p(y_n | \mathbf{x}_n, \mathbf{w}) = \begin{cases} \sigma(\mathbf{w}^T \mathbf{x}_n), & y_n = 1 \\ 1 - \sigma(\mathbf{w}^T \mathbf{x}_n), & y_n = 0 \end{cases} = [\sigma(\mathbf{w}^T \mathbf{x}_n)]^{y_n} [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]^{1-y_n}$$

Data likelihood for all training data:

$$L(\mathcal{D} | \mathbf{w}) = \prod_{n=1}^N p(y_n | \mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^N [\sigma(\mathbf{w}^T \mathbf{x}_n)]^{y_n} [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]^{1-y_n}$$

Log-likelihood for all training data:

$$l(\mathcal{D} | \mathbf{w}) = \sum_{n=1}^N \{y_n \log [\sigma(\mathbf{w}^T \mathbf{x}_n)] + (1 - y_n) \log [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\}$$

Cross-entropy error (negative log-likelihood):

$$\mathcal{E}(\mathbf{w}) = - \sum_{n=1}^N \{y_n \log [\sigma(\mathbf{w}^T \mathbf{x}_n)] + (1 - y_n) \log [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\}$$

Logistic Regression - Optimization

No closed-form solution that minimizes the cross-entropy function.

We use an approximate method, e.g. gradient descent, so we need to compute $\nabla \mathcal{E}(\mathbf{w})$.

Derivation of $\nabla \mathcal{E}(\mathbf{w}) = \frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}}$

$$\begin{aligned}\nabla \mathcal{E}(\mathbf{w}) &= - \sum_{n=1}^N \left\{ y_n [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)] \mathbf{x}_n - (1 - y_n) [1 - (1 - \sigma(\mathbf{w}^T \mathbf{x}_n))] \mathbf{x}_n \right\} \\ &= - \sum_{n=1}^N \left\{ y_n [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)] \mathbf{x}_n + (1 - y_n) \sigma(\mathbf{w}^T \mathbf{x}_n) \mathbf{x}_n \right\} \\ &= - \sum_{n=1}^N [y_n - y_n \sigma(\mathbf{w}^T \mathbf{x}_n) - \sigma(\mathbf{w}^T \mathbf{x}_n) + y_n \sigma(\mathbf{w}^T \mathbf{x}_n)] \mathbf{x}_n \\ &= \sum_{n=1}^N \underbrace{(\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n)}_{\text{error}} \mathbf{x}_n\end{aligned}$$

Gradient descent update: $\mathbf{w}_{k+1} := \mathbf{w}_k - \alpha(k) \nabla \mathcal{E}(\mathbf{w})$

Is the cross-entropy error a convex function?

Derivation of $\mathbf{H} = \frac{\partial^2 \mathcal{E}(\mathbf{w})}{\partial^2 \mathbf{w}} = \nabla \left((\nabla \mathcal{E}(\mathbf{w}))^T \right) = \nabla \left(\sum_{n=1}^N (\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n) \mathbf{x}_n^T \right)$

$$\begin{aligned}\mathbf{H} &= \frac{\partial}{\partial \mathbf{w}} \left[\sum_{n=1}^N (\sigma(\mathbf{w}^T \mathbf{x}_n) \cdot \mathbf{x}_n^T - y_n \mathbf{x}_n^T) \right] \\ &= \sum_{n=1}^N \frac{\partial}{\partial \mathbf{w}} [\sigma(\mathbf{w}^T \mathbf{x}_n)] \cdot \mathbf{x}_n^T \quad (\text{chain rule}) \\ &= \sum_{n=1}^N \underbrace{\sigma(\mathbf{w}^T \mathbf{x}_n)}_{\in [0,1]} \cdot \underbrace{(1 - \sigma(\mathbf{w}^T \mathbf{x}_n))}_{\in [0,1]} \cdot \underbrace{(\mathbf{x}_n \cdot \mathbf{x}_n^T)}_{\in \mathcal{R}^{D \times D}}\end{aligned}$$

For all $\mathbf{v} \in \mathbb{R}^D$, substituting $\mu_n = \sigma(\mathbf{w}^T \mathbf{x}_n) (1 - \sigma(\mathbf{w}^T \mathbf{x}_n)) \geq 0$, we have:

$$\mathbf{v}^T \mathbf{H} \mathbf{v} = \cdot \mathbf{v}^T \left(\sum_{n=1}^N \mu_n \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{v} = \sum_{n=1}^N (\mu_n \mathbf{x}_n^T \mathbf{v})^T (\mathbf{x}_n^T \mathbf{v}) = \sum_{n=1}^N \mu_n \|\mathbf{x}_n^T \mathbf{v}\|_2^2 \geq 0$$